

The new CMS DAQ system for LHC operation after 2014 (DAQ2)



CHEP2013: Computing in High Energy Physics 2013
14-18 Oct 2013
Amsterdam



Andre Holzner, University California at San Diego

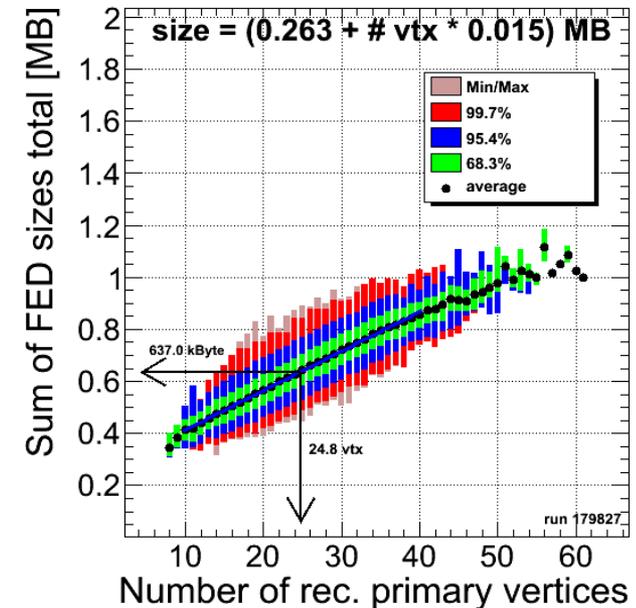
On behalf of the CMS collaboration

Overview

- DAQ2 Motivation
- Requirements
- Layout / Data path
- Frontend Readout Link
- Event builder core
- Performance considerations
- Infiniband
- File based filter farm and storage
- DAQ2 test setup and results
- Summary/Outlook

DAQ2 motivation

- Aging equipment:
 - Run1 DAQ uses some technologies which are disappearing
 - PCI-X cards, Myrinet
 - Almost all equipment reached the end of the 5 year lifecycle
- CMS detector upgrades
 - Some subsystems move to new front-end drivers
 - Some subsystems will add more channels
- LHC performance
 - Expect higher instantaneous luminosity after LS1
 - higher number of interactions per bunch crossing ('pileup')
 - larger event size, higher data rate
- Physics
 - Higher centre-of-mass energy and more pileup imply:
 - either raise trigger thresholds, or
 - more intelligent decisions at Higher Level Trigger
 - requires more CPU power



DAQ2 requirements

Requirement	DAQ1	DAQ2
Readout rate	100 kHz	
Front end drivers (FEDs)	640: 1 ~ 2 kByte	640: 1 ~ 2 kByte ~50: 2 ~ 8 kByte
Total readout bandwidth	100 GByte/s	200 GByte/s
Interface to FEDs ¹⁾	SLink64	Slink64/Slink Express
Coupling Event builder software/HLT software ²⁾	no requirement	decoupled
Lossless event building	<input type="checkbox"/>	
HLT capacity	extendable	
High availability/ fault tolerance ³⁾	<input type="checkbox"/>	
Cloud facility for offline processing ⁴⁾	originally not required	<input type="checkbox"/>
Subdetector local runs	<input type="checkbox"/>	

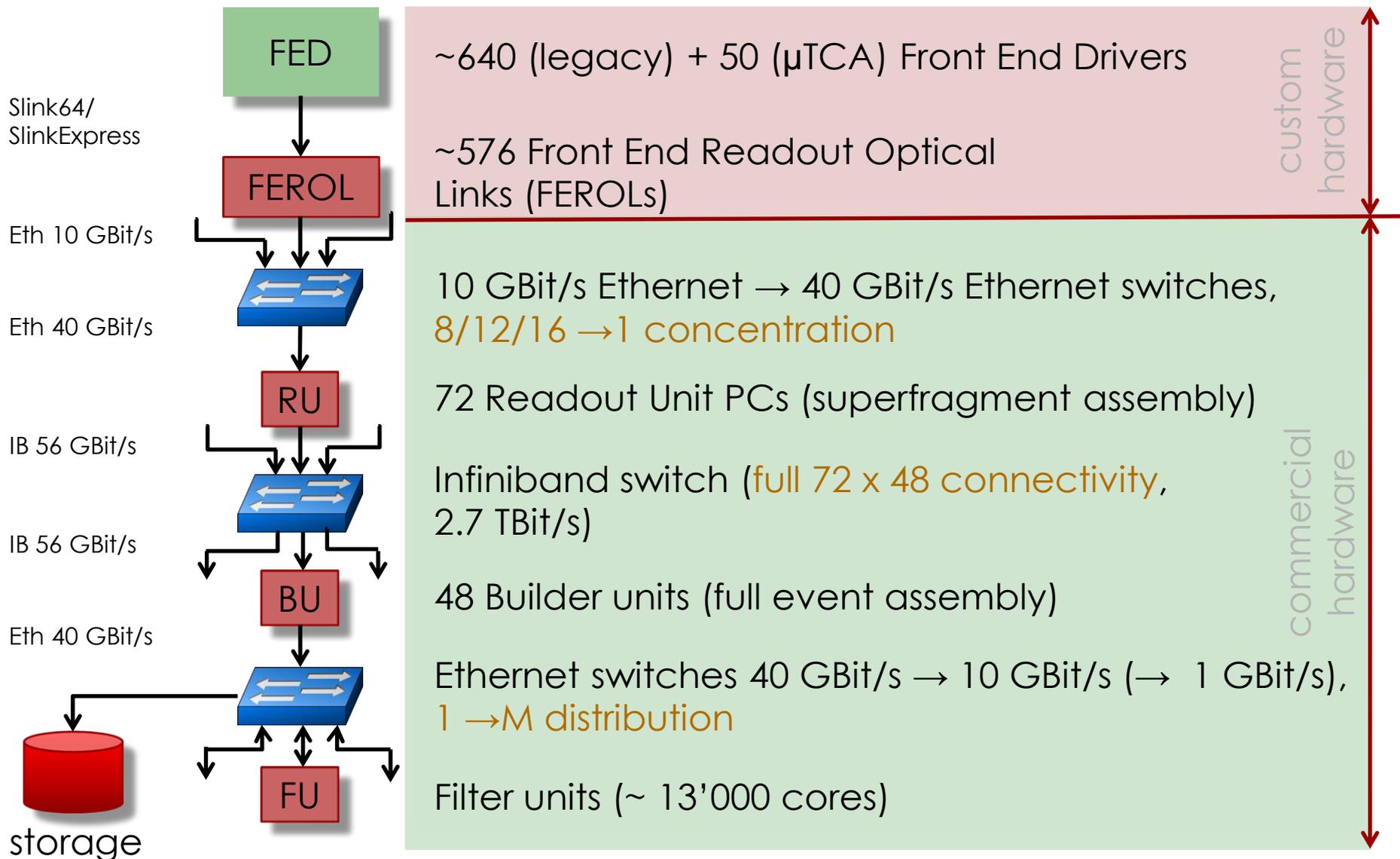
See talks of 1) P. Žejdl

2) R. Mommsen

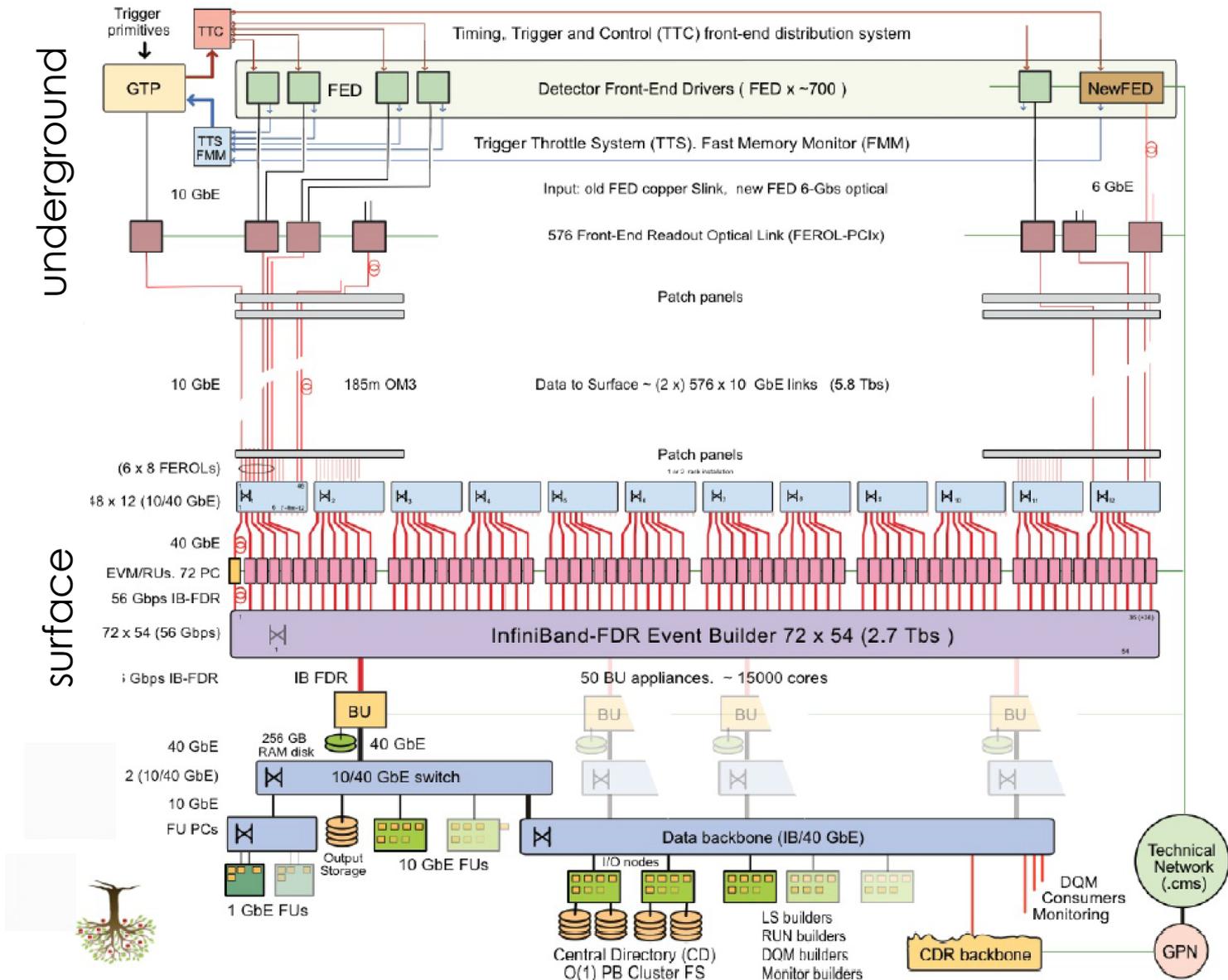
3) H.Sakulin

4) J.A.Coarasa

DAQ2 data path



DAQ2 layout



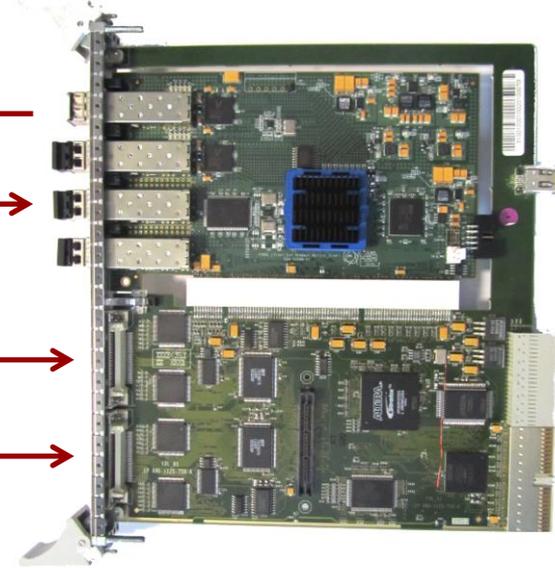
FrontEnd Readout Link (FEROL)

- Replace Myrinet card (upper half) by a new custom card
- PCI-X interface to legacy slink receiver card (lower half)
- 10 GBit/s Ethernet output to central event builder
 - Restricted TCP/IP protocol engine inside the FPGA
- Additional optical links (inputs) for future μ TCA based Front End Drivers (6-10 GBit/s; custom, simple point to point protocol)
- Allows to use industry standard 10 GBit/s transceivers, cables and switches/routers
 - Only commercially available hardware further downstream

10 GBit/s Ethernet

Slink Express
from μ TCA
FEDs

Slink64
from FEDs

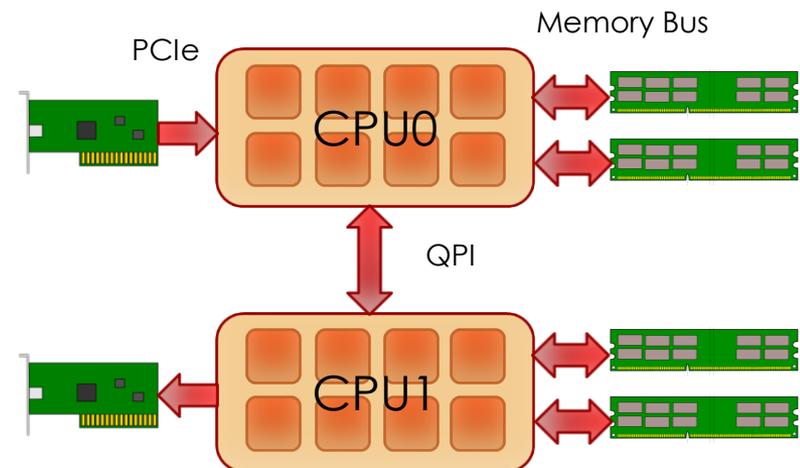


see P. Žejdl's talk
for more details

Performance considerations

- Number of DAQ2 elements is an order of magnitude smaller than for DAQ1
- Consequently, bandwidth per PC is an order of magnitude higher
- CPU frequency did not increase since DAQ1 but number of cores did
- Need to pay attention to performance tuning
 - TCP socket buffers
 - Interrupt affinities
 - Non-uniform memory access

	DAQ1	DAQ2
# readout units (RU)	640	48
RU max. bandwidth	3 Gbit/s	40 Gbit/s
# builder units (BU)	>1000	72
BU max. bandwidth	2 Gbit/s	56 Gbit/s



Infiniband

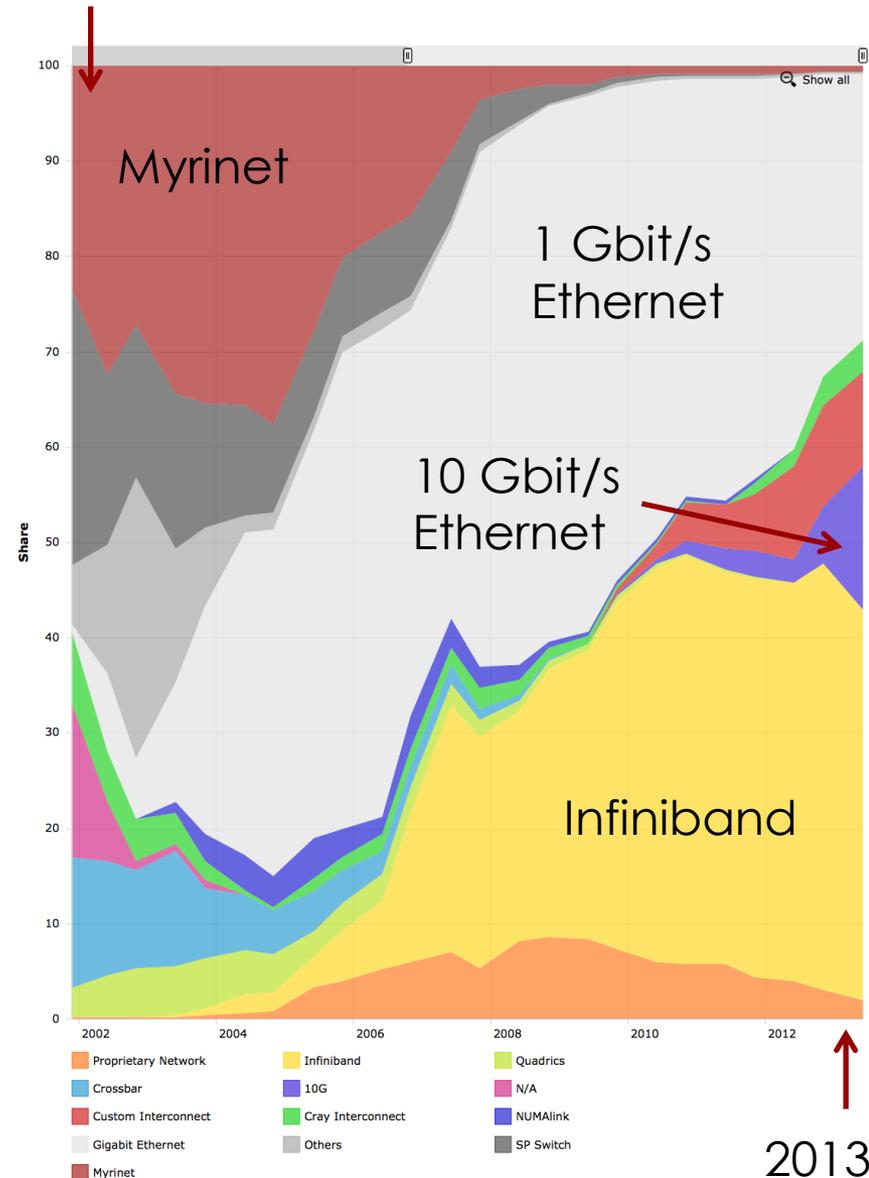
DAQ1 TDR (2002)

Advantages:

- Designed as a High Performance Computing interconnect over short distances (within datacenters)
- Protocol is implemented in the network card silicon → low CPU load
- 56 GBit/s per link (copper or optical)
- Native support for Remote Direct Memory Access (RDMA)
- No copying of bulk data between user space and kernel ('true zero-copy')
- affordable

Disadvantages:

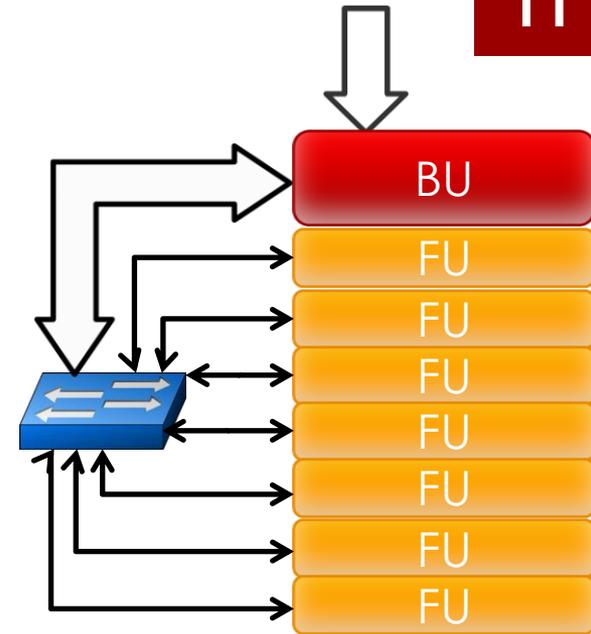
- Less widely known, API significantly differs from BSD sockets for TCP/IP
- Fewer vendors than Ethernet
- Niche market



Top500.org share by Interconnect family

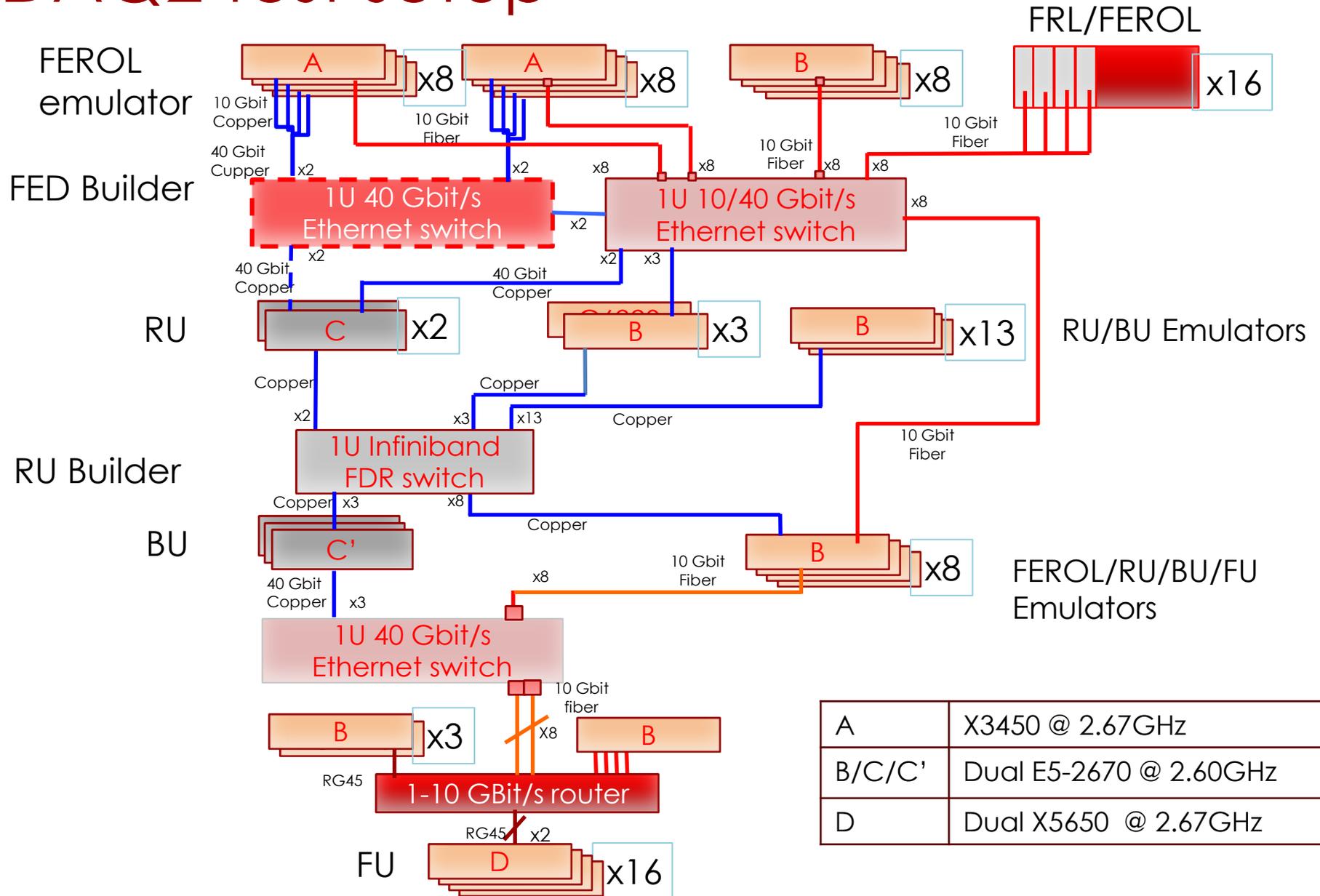
File based filter farm and Storage

- In DAQ1, high level trigger process was running inside a DAQ application
 - introduces dependencies between online (DAQ) and offline (event selection) software which have different release cycles, compilers, state machines etc.
- Decoupling these needs a common, simple interface
 - files (no special common code required to write and read them)
- Builder unit stores events in files in a RAM disk
 - Builder Unit acts as a NFS server, exports event files to Filter Unit PCs
 - Baseline: 2 Gbyte/s bandwidth
 - 'Local' within a rack
- Filter units write out selected events (~ 1 in 100) back to a global (CMS DAQ wide) filesystem (e.g. Lustre) for transfer to Tier0 computing centre



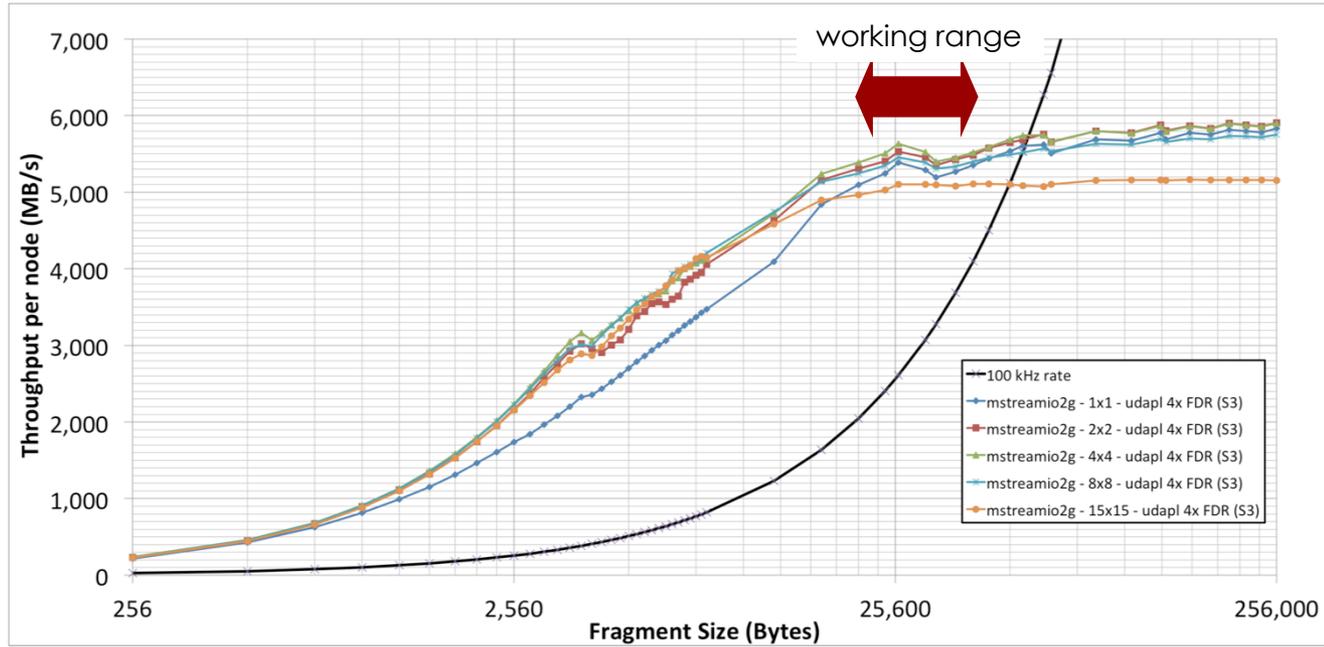
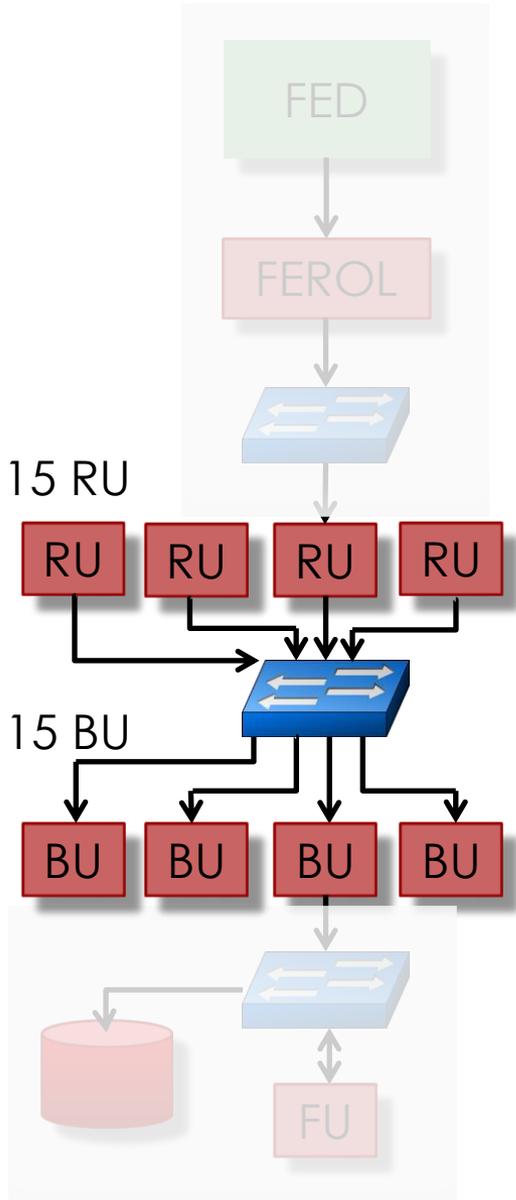
see R. Mommsen's talk
for more details

DAQ2 test setup

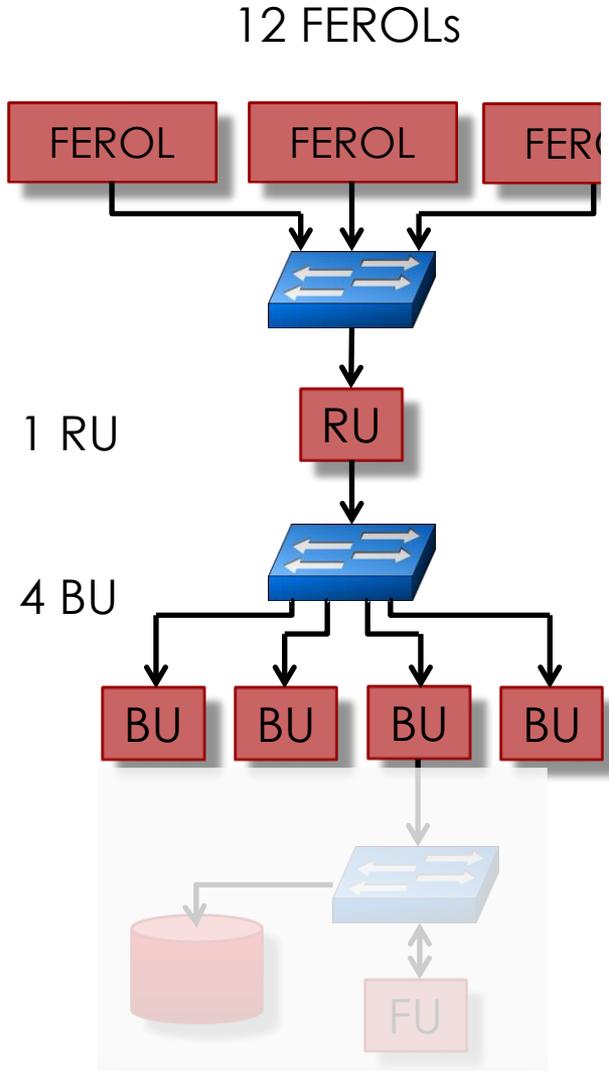


A	X3450 @ 2.67GHz
B/C/C'	Dual E5-2670 @ 2.60GHz
D	Dual X5650 @ 2.67GHz

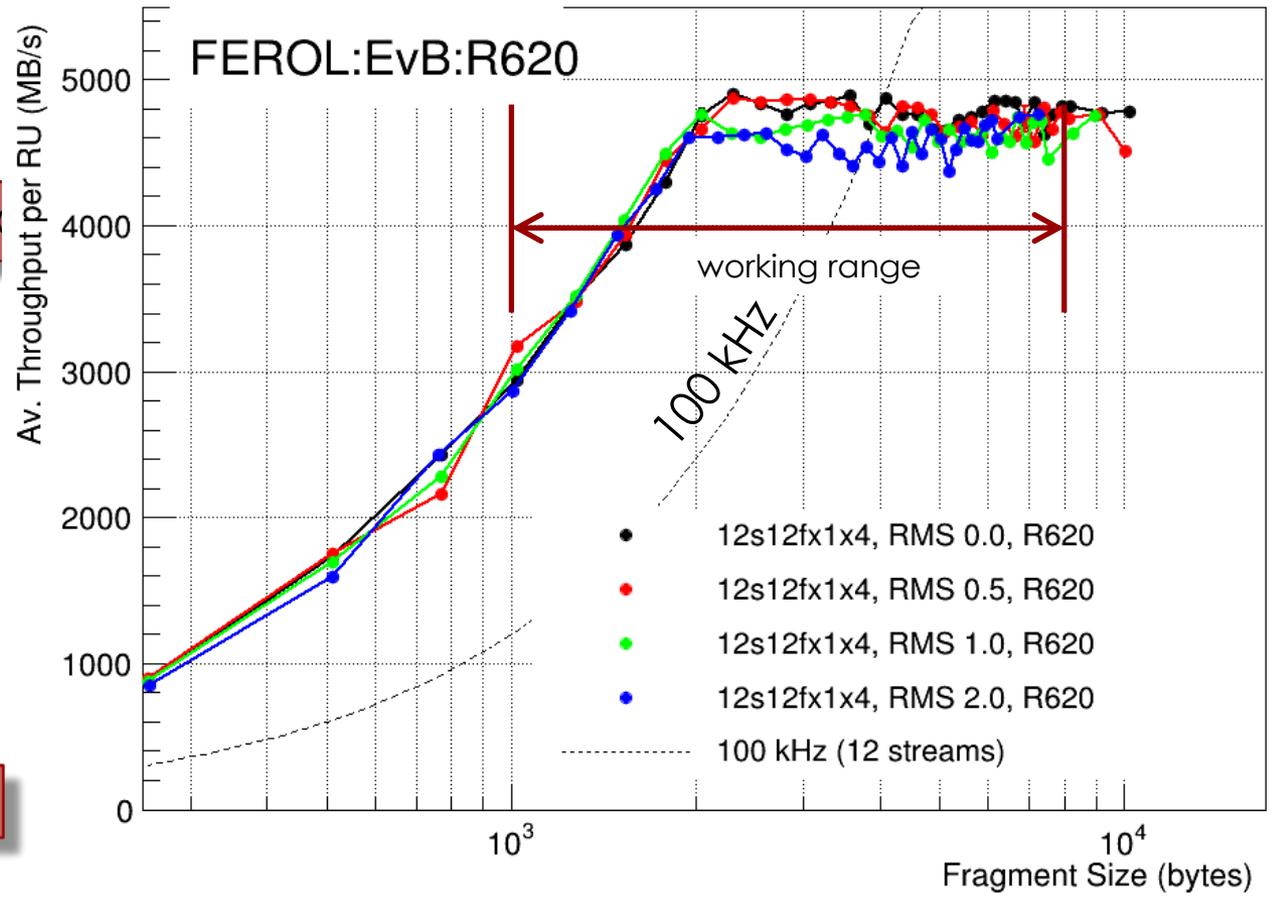
InfiniBand Measurements



Test setup results

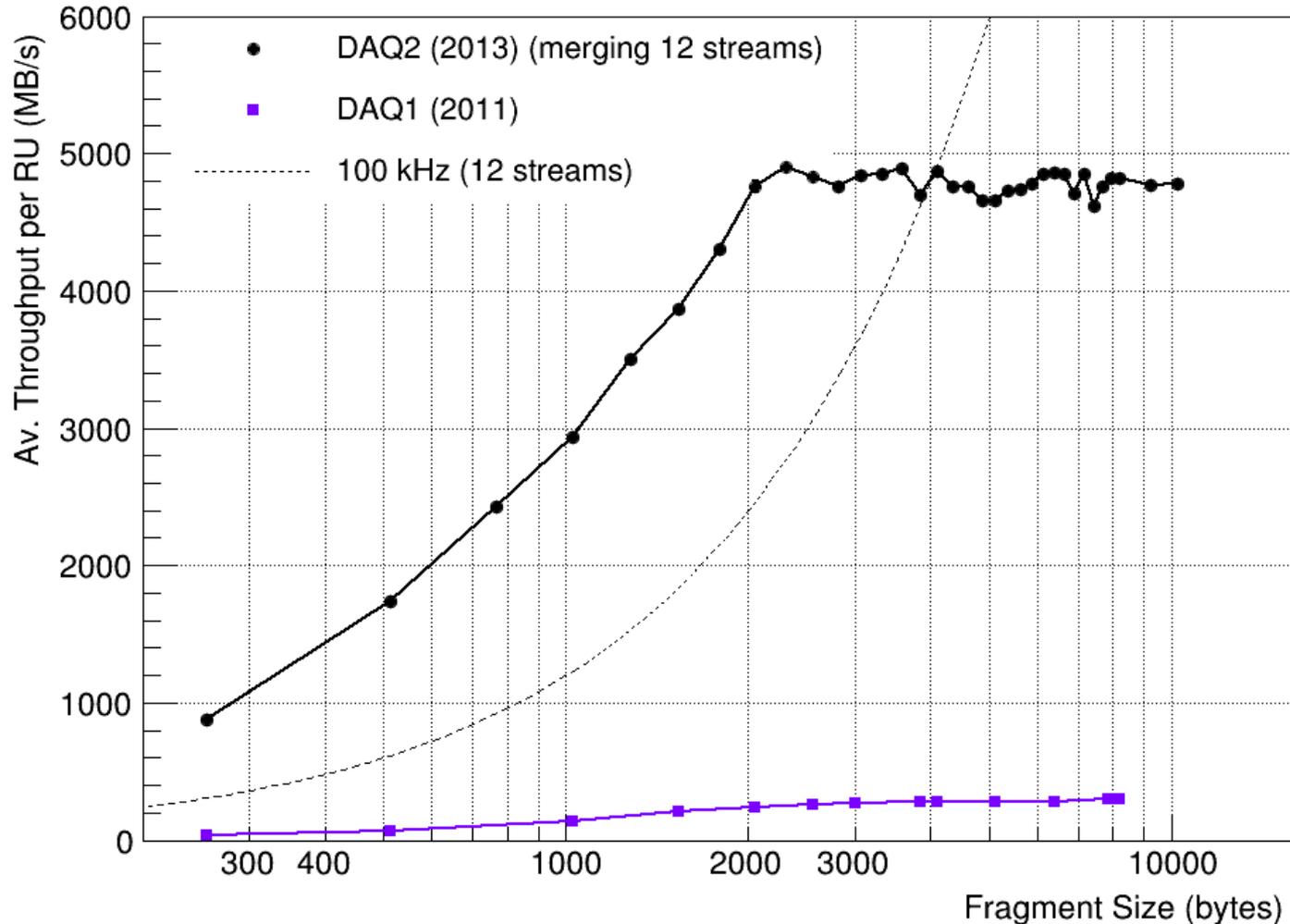


Throughput vs. Fragment Size



Test setup: DAQ1 vs. DAQ2

Throughput vs. Fragment Size



Comparison of throughput per Readout Unit

Summary / Outlook

- CMS has designed a central data acquisition system for post-LS1 data taking
 - replacing outdated standards by modern technology
 - ~ twice the event building capacity than DAQ system for Run1
 - accomodating a large dynamic range of up to 8 kByte fragments, flexible configuration
- Increase in networking bandwidth was faster than increase in event sizes
 - Number of event builder PCs reduced by a factor ~10
 - Each PC handles a factor ~10 more bandwidth
 - Requires performance related fine-tuning
- Performed various performance tests with a small scale demonstrator
- First installation activities for DAQ2 have started already
 - Full deployment foreseen for mid 2014
- Looking forward to recording physics data after the Long Shutdown 1 !